

Un index de jointure pour les entrepôts de données XML

Hadj Mahboubi, Kamel Aouiche, Jérôme Darmont

ERIC, Université Lumière Lyon 2
5 avenue Pierre Mendès-France
69676 Bron Cedex
{ hmahboubi | kaouiche | jdarmont } @eric.univ-lyon2.fr

Résumé. Les entrepôts de données XML proposent une base intéressante pour les applications décisionnelles qui exploitent des données hétérogènes et provenant de sources multiples. Cependant, les performances des SGBD natifs XML étant actuellement limitées, il est nécessaire de trouver des moyens de les optimiser. Dans cet article, nous proposons un nouvel index spécifiquement adapté à l'architecture multidimensionnelle des entrepôts de données XML, qui élimine le coût des jointures tout en préservant l'information contenue dans l'entrepôt initial. Une étude théorique et des résultats expérimentaux démontrent l'efficacité de notre index, même lorsque les requêtes sont complexes.

1 Introduction

Les technologies entrant en compte dans les processus décisionnels, comme les entrepôts de données (*data warehouses*), l'analyse multidimensionnelle en ligne (*On-Line Analysis Process* ou OLAP) et la fouille de données (*data mining*), sont désormais très efficaces pour traiter des données simples, numériques ou symboliques. Cependant, les données exploitées dans le cadre des processus décisionnels sont de plus en plus complexes. L'avènement du Web et la profusion de données multimédia ont en grande partie contribué à l'émergence de cette nouvelle sorte de données. Dans ce contexte, le langage XML peut grandement aider à l'intégration et à l'entreposage de ces données. C'est pourquoi nous nous intéressons aux travaux émergents sur les entrepôts de données XML (Wolfgang et al., 2003; Baril et Bellahsène, 2003; Pokorný, 2001; Golfarelli et al., 2001). Cependant, les requêtes décisionnelles exprimées en XML sont généralement complexes du fait qu'elles impliquent de nombreuses jointures et agrégations. Par ailleurs, les systèmes de gestion de bases de données (SGBD) natifs XML présentent actuellement des performances médiocres quand les volumes de données sont importants ou que les requêtes sont complexes. Il est donc crucial lors de la construction d'un entrepôt de données XML de garantir la performance des requêtes XQuery qui l'exploiteront.

Plusieurs études traitent de l'indexation des données XML (K.Gupta et al.; Yeh et Gardarin, 2001; Chung et al., 2002). Ces index optimisent principalement des requêtes exprimées en expressions de chemin. Or, dans le contexte des entrepôts de données XML, les requêtes sont complexes et comportent plusieurs expressions de chemin. De plus, ces index opèrent sur un seul document et ne prennent pas en compte d'éventuelles jointures, qui sont courantes dans les requêtes décisionnelles. À notre connaissance, seul l'index Fabric (Cooper et al., 2001) permet actuellement gérer plusieurs documents XML. Cependant, cet index ne prend pas en compte

Un index de jointure pour les entrepôts de données XML

les relations qui peuvent exister entre ces documents (les prédicats de jointure, notamment) et n'est donc pas non plus adapté à nos besoins.

C'est pourquoi nous proposons une structure d'index spécifiquement adaptée aux données multidimensionnelles d'un entrepôt XML, c'est-à-dire une structure capable de maintenir des données de plusieurs documents XML à la fois tout en préservant l'information contenue dans ces données et leur modélisation en étoile. Notre structure d'index, que nous qualifions d'index de jointure, permet également d'assurer un meilleur traitement des requêtes décisionnelles exprimées en XQuery en éliminant les coûts de jointure. Afin de valider cette proposition, nous avons mené une étude théorique ainsi que des expérimentations sur un entrepôt de données XML réel.

Le reste de cet article est organisé comme suit. Nous présentons dans la section 2 le contexte de notre étude ainsi que notre structure d'index. Les études théorique et expérimentale que nous avons menées pour tester la validité de notre index sont présentées dans la section 3. Finalement, nous concluons et discutons nos perspectives de recherche dans la section 4.

2 Indexation des entrepôts de données XML

2.1 Contexte

Afin d'appliquer notre démarche d'optimisation des performances, nous avons sélectionné l'architecture d'entrepôt de données XCube (Wolfgang et al., 2003), qui propose une modélisation en étoile de données stockées dans des documents XML. Ces documents permettent de représenter respectivement le schéma et les métadonnées (*Schema.xml*), les dimensions (*Dimensions.xml*) et la table de faits (*TableFacts.xml*) de l'entrepôt XML. La figure 1 (a) et 1 (b) représentent les deux derniers documents sous forme de graphes XML.

Pour l'interrogation de l'entrepôt, nous avons adopté le langage XQuery car il permet de formuler des requêtes complexes. La requête **Q** donne un exemple d'une requête décisionnelle exprimée dans ce langage. Cette requête retourne la moyenne des quantités vendues aux clients de Lyon, avec regroupement par nom et codes postal. Elle réalise une jointure entre le document *Dimensions.xml* et *TableFacts.xml*. Notons que XQuery ne permet pas normalement de faire des opérations de groupement *Group by* multiples. Dans notre implémentation, nous avons donc étendu l'interface d'interrogation du SGBD natif XML eXist (Meier, 2002) de manière à ce qu'il puisse traiter ce type de requêtes.

```
Q : for $a in //dimensionData/classification/Level[@node='customers']/node, $x in //CubeFacts/cube/Cell
where $b/attribute[@name='cust_city',@value='Lyon'] and $x/dimension //@node=$a/@id and $x/dimension
//@id='customers'
group by(@cust_first_name,@cust_postal_code) return sum(quantity)
```

2.2 Structure de notre index de jointure

Notre index doit permettre de conserver les relations entre les dimensions et les mesures dans un fait. Il présente donc une structure similaire à celle du document *TableFacts.xml*, à l'exception de l'élément *attribute*.

Sa structure est présentée dans la figure 1 (c). Les étiquettes qui commencent par le caractère @ représentent les attributs et les autres représentent les éléments. Chaque cellule est

identifiée par des dimensions et un ou plusieurs faits. Un fait (élément *Fact*) possède deux attributs, *@id* et *@value*, qui indiquent respectivement son nom et sa valeur. Chaque dimension (élément *dimension*) est identifiée par deux attributs : *@id* qui donne le nom de la dimension et *@node* qui donne la valeur de l'identifiant de la dimension. De plus, l'élément dimension possède un certain nombre d'éléments fils (éléments *attribute*). Ces éléments sont insérés pour stocker les noms et les valeurs des attributs de chaque dimension. Ils sont obtenus depuis le document *Dimensions.xml*. Un élément *attribute* est caractérisé par deux attributs, *@nom* et *@value*, qui indiquent respectivement le nom et la valeur de chaque attribut.

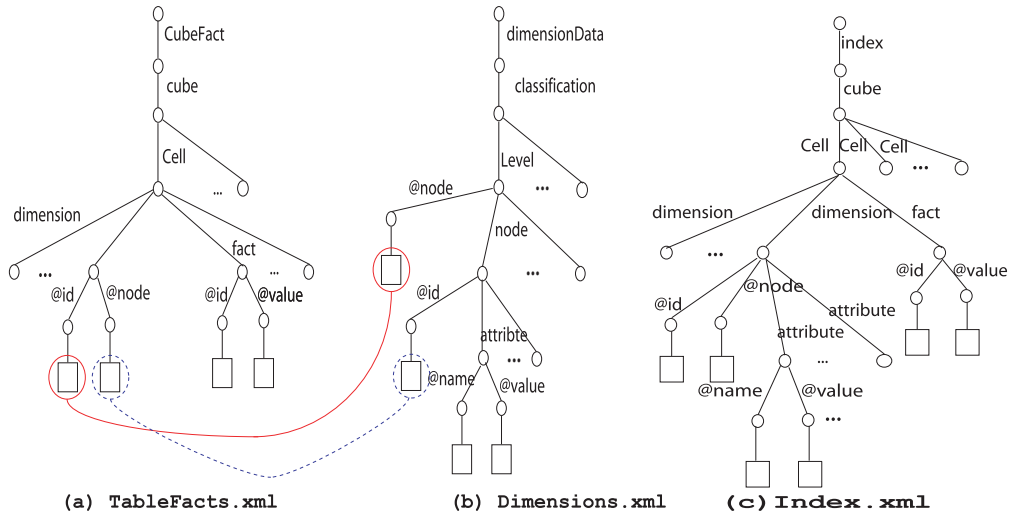


FIG. 1 – Structure en arbre des documents (a) *Dimensions.xml*, (b) *TableFacts.xml* et (c) *Index.xml*

La migration des données des documents *Dimensions.xml* et *TableFacts.xml* vers la structure d'index, et le fait de stocker dans une même cellule les faits, les dimensions et leurs attributs, nous permet d'éliminer les opérations de jointure. Toutes les informations nécessaires pour la jointure sont en effet stockées dans la même cellule.

3 Validation

3.1 Étude théorique

Une requête type définie sur un entrepôt de données XML modélisé selon la spécification XCube réalise plusieurs jointures entre les faits stockés dans *TableFacts.xml* et les dimensions de *Dimensions.xml*. Il faut alors vérifier les contraintes suivantes : $TableFacts.@id = Dimensions.@node$ et $TableFacts.@node = Dimensions.@id$. La première égalité vérifie que la dimension composant une cellule est bel et bien la dimension exprimée dans la requête. La seconde vérifie que le nœud d'une dimension (équivalent d'une clé primaire) correspond

Un index de jointure pour les entrepôts de données XML

(peut être joint) au nœud, de la même dimension, défini dans une cellule (équivalent d'une clé étrangère de la table de faits).

L'exécution d'une requête sans utilisation de notre index peut se dérouler comme suit. Pour chaque dimension définie par $@node = 'nom\ de\ la\ dimension'$, les identifiants $@id$ vérifiant la clause *Where* sont recherchés. Le document *Dimensions.xml* est parcouru en profondeur, jusqu'au nœud *Level*. Les fils *node* du nœud *Level* sont ensuite parcourus en largeur jusqu'à trouver le nœud dont la valeur de $@node$ est égale au nom de la dimension spécifié dans la requête. Le coût de ce parcours est égal au nombre de nœuds *Level* du document *Dimensions.xml*, c'est-à-dire le nombre de dimensions dans le schéma, dénoté $|dimension|$. Si plusieurs dimensions sont définies, le parcours donne alors lieu à autant de nœuds que de dimensions. En revanche, le coût de ce parcours reste invariant car tous les nœuds *Level* sont parcourus. Pour chaque nœud trouvé, ses fils sont parcourus en profondeur jusqu'à trouver la liste des $@id$ vérifiant les conditions $@name = 'nom\ de\ l'attribut'$ et $@value = 'valeur\ de\ l'attribut'$. Le coût de ce parcours est égal au nombre de fils *attribute*. En résumé, le coût de traitement des dimensions est égal à $|a_i| * |d_i|$, où $|a_i|$ désigne le nombre d'attributs de chaque dimension et $|d_i|$ le nombre d'éléments *node*, c'est-à-dire le nombre de fils de chaque dimension. Pour réaliser une jointure entre les dimensions du document *Dimensions.xml* et les faits du document *TableFacts.xml*, les $@id$ retrouvés dans le traitement des dimensions sont recherchés dans les faits. Le document *TableFacts.xml* est alors parcouru en profondeur jusqu'au niveau *Cell*. Le coût de ce parcours est égal à 2. Les cellules sont ensuite parcourues en largeur afin de trouver les dimensions dont le fils $@id$ est égal à $@node$ de *Dimensions.xml* et $@node$ est égal à $@id$ de *Dimensions.xml*. En résumé, le coût de traitement du document *TableFacts.xml* est $|Cell|$, où $|Cell|$ est le nombre de cellules du document *TableFacts.xml*. Le coût de traitement d'une requête est donnée par la formule $E_{sans-index} = ((|Cell|) * |Dimension|) * (|Dimension| + (|d_i| * |a_i|))$.

L'exécution d'une requête, avec utilisation de notre index (stocké dans le document *Index.xml*), peut se dérouler comme suit. Pour chaque dimension définie par $@node = 'nom\ de\ la\ dimension'$, les identifiants $@id$ vérifiant la clause *Where* sont recherchés. Le document *Index.xml* est parcouru en profondeur, jusqu'au nœud *Cell*. Le coût de ce parcours est égal au nombre de cellules dans le document *Index.xml*. Les fils *dimension* du nœud *Cell* sont ensuite parcourus en largeur jusqu'à trouver le nœud dont la valeur de $@id$ est égale au nom de la dimension spécifié dans la requête. Le coût de ce parcours est égal au nombre de nœuds *dimension* du document *Index.xml*, c'est-à-dire le nombre de dimensions dans le schéma de l'entrepôt de données : $|dimension|$. Pour chaque nœud trouvé, ses fils sont parcourus en profondeur jusqu'à trouver le nœud *attribute* vérifiant les conditions $@name = 'nom\ de\ l'attribut'$ et $@value = 'valeur\ de\ l'attribut'$. Le coût de ce parcours est égal au nombre de fils *attribute* $|a_i|$. En résumé, le coût de traitement d'une requête qui exploite notre structure d'index est donné par la formule $E_{index} = |Cell| * (|Dimension| + |a_i|)$.

La figure 2 (a) représente la variation des coûts $E_{sans-index}$ et E_{index} en fonction du nombre de cellules. Nous constatons que l'utilisation de notre index permet un gain de facteur 14000 en moyenne.

3.2 Expérimentations

En complément de notre étude théorique, nous avons effectué des expérimentations afin de tester l'efficacité de notre proposition de structure d'index. Nous avons généré un entrepôt de données XCube implanté au sein du SGBD XML natif eXist. Construit à partir d'un entrepôt de

données relationnel test classique, cet entrepôt est constitué d'une table de faits *sales* (4,92 Mo) et de cinq dimensions *channels*, *promotions*, *customers*, *products* et *times* (3,77 Mo). Nous avons effectué nos tests sur une machine dotée d'un processeur Intel Pentium 4 GHz avec 1 GB de mémoire et un disque dur IDE. Nous avons exécuté la requête décisionnelle XQuery de la section 2.1 avec et sans utilisation de notre index et en faisant varier la taille de l'entrepôt. La figure 2 (b) présente les résultats obtenus exprimés en temps d'exécution par rapport au nombre de cellules (faits) dans le document *TableFacts.xml*.

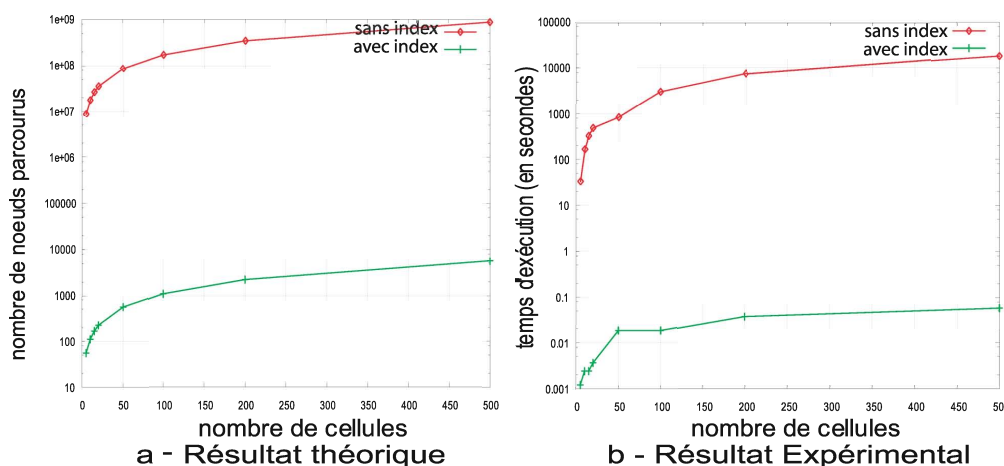


FIG. 2 – Résultats

La figure 2 (b) montre qu'avec l'utilisation de notre index, nous obtenons des temps de traitement en moyenne 25669 fois inférieurs à ceux obtenus sans utiliser notre index. De plus, cette figure est semblable à celle qui représente nos estimations théoriques (figure 2 (a)). Nous avons également utilisé notre structure d'index pour traiter la requête sur la totalité des cellules du documents *TableFacts.xml*. Nous avons obtenu un temps d'exécution de moins de deux secondes, alors que le système s'est avéré incapable de répondre à la requête dans un temps raisonnable lorsque nous n'utilisons pas notre index.

4 Conclusion et perspectives

Nous avons présenté dans cet article un nouvel index de jointure spécifiquement adapté aux entrepôts de données XML. Cette structure de données permet d'optimiser les temps d'accès à plusieurs documents XML en éliminant le coût de jointure, tout en préservant l'information contenue dans l'entrepôt initial. Une étude de complexité et des expérimentations nous ont permis de démontrer l'efficacité de notre index même lorsque les requêtes sont complexes et la taille de l'entrepôt volumineuse. Ces expérimentations nous ont par ailleurs amenés à étendre la syntaxe des requêtes XQuery afin qu'elles puissent supporter des clauses de regroupement (*Group by*) multiples.

Les perspectives de ce travail se situent dans le cadre du développement des SGBD natifs XML, qui vise à les doter des mêmes fonctionnalités et performances que les SGBD relationnels.

Premièrement, notre index pourrait être directement intégré au sein d'un SGBD natif XML. Il serait également indispensable de mettre au point une stratégie de maintenance incrémentale de la structure de données de notre index lorsque les données sources sont mises à jour. Par ailleurs, notre mécanisme de réécriture de requêtes pourrait également être directement intégré, notamment au niveau de la modélisation des vues. Finalement, les modifications que nous avons apportées à la clause de regroupement *Group by* de XQuery pourraient aussi faire objet d'une proposition d'extension de la syntaxe de ce langage.

Références

- Baril, X. et Z. Bellahsène (2003). *Designing and Managing an XML Warehouse*, pp. 455–473. XML Data Management: Native XML and XML-enabled Database Systems. Addison Wesley.
- Chung, C., J. Min, et K. Shim (2002). APEX: An adaptive path index for XML data. In *ACM SIGMOD International Conference on Management of Data (SIGMOD 2002)*, Madison, USA, pp. 121–132.
- Cooper, B., N. Sample, M. J. Franklin, G. R. Hjaltason, et M. Shadmon (2001). A fast index for semistructured data. In *27th International Conference on Very Large Data Bases (VLDB 2001)*, Roma, Italy, pp. 341–350.
- Golfarelli, M., S. Rizzi, et B. Vrdoljak (2001). Data warehouse design from XML sources. In *4th ACM international workshop on Data warehousing and OLAP (DOLAP 2001)*, Atlanta, USA.
- K.Gupta, R., G.Shuqiao, et Y.Zhen. A report on XML data indexing techniques. Technical report, National University of Singapore.
- Meier, W. (October 2002). eXist: An open source native XML database. In *Akmal B. Chaudri, Mario Jeckle, Erhard Rahm, Rainer Unland (Eds.): Web, Web-Services, and Database Systems. NODe 2002 Web- and Database-Related Workshops*, Erfurt, Germany.
- Pokorný, J. (2001). Modelling stars using XML. In *4th ACM international workshop on Data warehousing and OLAP (DOLAP 2001)*, Atlanta, USA, pp. 24–31.
- Wolfgang, H. Andreas, et B. Harde (2003). XCube: XML for data warehouses. In *6th ACM International Workshop on Data warehousing and OLAP (DOLAP 2003)*, New Orleans, USA, pp. 33–40.
- Yeh, L. et G. Gardarin (2001). Indexing XML objects with ordered schema trees. In *20th Conference on Bases de Données Avancées (BDA 2001)*, Montpellier, France, pp. 361–370.

Summary

XML data warehouses form an interesting basis for decision-support applications that exploit heterogeneous data from multiple sources. However, XML-native database systems currently bear limited performances and it is necessary to research ways to optimize them. In this paper, we propose a new index that is specifically adapted to the multidimensional architecture of XML warehouses and eliminates join operations, while preserving the information

contained in the original warehouse. A theoretical study and experimental results demonstrate the efficiency of our index, even when queries are complex.